

VISVESVARAYA TECHNOLOGICAL UNIVERSITY
“Jnana Sangama”, Belgavi-590018



A Project Report
On
**“NAMED ENTITY RECOGNITION IN SOCIAL MEDIA
USING MACHINE LEARNING”**

*Submitted in the partial fulfillment for the award of the Bachelor of Engineering degree in
Computer Science and Engineering*

Submitted By

| | |
|-------------------|------------|
| AISHWARYA M S | 4AD21CS004 |
| ESHWAR A M | 4AD21CS021 |
| KUSHAL R | 4AD21CS039 |
| LAKSHMI PRIYA H P | 4AD21CS040 |

Under the guidance of

Mrs. Sushma V
Assistant Professor
Department of Computer Science & Engineering
ATME College of Engineering



A T M E
College of Engineering

Department of Computer Science & Engineering
ATME College of Engineering
13th Kilometer, Mysore – Kanakapura - Bangalore Road,
Mysore-570028
2024-2025

VISVESVARAYA TECHNOLOGICAL UNIVERSITY
“Jnana Sangama”, Belgavi-590018



Department of Computer Science and Engineering



A T M E
College of Engineering

CERTIFICATE

This is to certify that the Project Work entitled “**NAMED ENTITY RECOGNITION IN SOCIAL MEDIA USING MACHINE LEARNING**” is the Bonafide work carried out by **Aishwarya M S (4AD21CS004)**, **Eshwar A M (4AD21CS021)**, **Kushal R (4AD21CS039)** and **Lakshmi Priya H P (4AD21CS040)** in partial fulfillment for the award of degree of Bachelor of Engineering in Computer Science and Engineering from Visvesvaraya Technological University, Belagavi during the year 2024-2025. The report has been approved and satisfies the academic requirement with respect to Project Work prescribed for Bachelor of Engineering degree.

Signature of Guide

Mrs. Sushma V

Assistant Professor

Signature of HOD

Dr. Puttegowda D

Professor & Head

Signature of Principal

Dr. L Basavaraj

Principal

External Viva

Name of Examiners

1.....

2.....

Signature with Date

.....

.....

ACKNOWLEDGEMENT

The successful completion of our work would be incomplete without the mention of the names of the people who have made it possible. We are indebted to several individuals who have helped us to complete our project.

We are thankful to **Dr. L Basavaraj, Principal**, ATME College of Engineering for having granted us permission and extended full use of the college facilities to carry out our project successfully.

We express our profound gratitude to **Dr. Puttegowda D**, Professor & Head, Department of Computer Science and Engineering for his consistent co-operation and support.

At the outset we express our profound gratitude to our guide **Mrs. Sushma V**, Assistant Professor, Department of Computer Science and Engineering for her consistent co-operation and support.

We are greatly indebted to our project coordinators **Dr. Sunitha Patel M S** and **Dr. Shilpa B L**, Associate Professors, Department of Computer Science and Engineering for her timely inquiries into the progress of the project

We also thank all teaching and non-teaching staff members of the department for their valuable assistance throughout our academic tenure.

Lastly, we would like to thank our family and friends for their cooperation and support for successful completion of our project.

| | |
|--------------------------|---------------------|
| Aishwarya M S | (4AD21CS004) |
| Eshwar A M | (4AD21CS021) |
| Kushal R | (4AD21CS039) |
| Lakshmi Priya H P | (4AD21CS040) |

ABSTRACT

The growing prevalence of social media platforms has led to an overwhelming amount of unstructured textual data, making it challenging to extract meaningful information. This project introduces a system for Named Entity Recognition (NER) on social media text using machine learning techniques, specifically leveraging a fine-tuned DistilBERT model. The system fetches posts from Reddit based on user-specified keywords and subreddits, processes the content, and identifies key entities such as names, organizations, and locations.

Using token classification with an aggregation strategy, the model analyzes noisy and informal text often prevalent in social media. Extracted entities are displayed in interactive visualizations, including frequency charts and entity type distributions, for better interpretability. The project highlights the challenges of working with noisy social media data and provides a robust, scalable solution for text analytics. By automating entity extraction, this system demonstrates the potential of machine learning in streamlining information retrieval from vast online datasets, with applications in market research, sentiment analysis, and digital investigations.

TABLE OF CONTENTS

| Chapter No. | Chapter Name | Page No. |
|--------------------|--|-----------------|
| | ACKNOWLEDGEMENT | i |
| | ABSTRACT | ii |
| | TABLE OF CONTENTS | iii-iv |
| | LIST OF FIGURES | v |
| Chapter 1 | Introduction | 1-2 |
| 1.1 | Overview | 1 |
| 1.2 | Existing System | 1 |
| 1.3 | Problem Statement | 1 |
| 1.4 | Proposed System | 1 |
| 1.5 | Advantages Over Current System | 2 |
| Chapter 2 | Literature Survey | 3-9 |
| 2.1 | Survey Papers | 3-8 |
| 2.2 | Related Work | 8-9 |
| Chapter 3 | System Requirements and Specification | 10-15 |
| 3.1 | Functional Requirements | 10 |
| 3.2 | Non-Functional Requirements | 11-15 |
| Chapter 4 | System Analysis | 16-17 |
| Chapter 5 | Methodology | 18-25 |
| 5.1 | System Architecture | 18-22 |
| 5.2 | Proposed System | 23-25 |
| Chapter 6 | Implementation | 26-34 |
| 6.1 | Stages of Implementation | 27-30 |
| 6.2 | Code Snippets | 31-36 |

| | | |
|------------------|--|--------------|
| Chapter 7 | System Testing | 37-37 |
| 7.1 | Testing Strategy | 37 |
| 7.2 | Unit Testing | 38 |
| 7.3 | Integration Testing | 39 |
| 7.4 | Validation Testing | 39 |
| 7.5 | Output Testing | 40 |
| 7.6 | System Testing | 40 |
| 7.7 | Performance and Scalability Testing | 41 |
| Chapter 8 | Results | 42-45 |
| 8.1 | Comparative Study | 42-45 |
| 8.2 | Snapshots | 46-47 |
| | Conclusion | 48 |
| | Future Scope | 49 |
| | References | 50-51 |

Chapter 1

INTRODUCTION

1.1 Overview

With the rise of social media platforms, a massive amount of textual data is generated daily. Extracting meaningful information from this unstructured data is vital for various applications, including sentiment analysis, trend monitoring, and targeted marketing. Named Entity Recognition (NER) is one such process that identifies and classifies entities in text. However, the informal and noisy nature of social media text poses challenges for traditional NER systems. This project addresses these challenges using BERT-based models.

1.2 Existing System

Traditional NER approaches, including rule-based systems and older machine learning models, are ill-equipped to handle social media data due to its informal syntax, abbreviations, and frequent misspellings. These systems lack the contextual understanding required to identify entities accurately.

1.3 Problem Statement

The primary objective of this project is to develop an NER system capable of accurately identifying named entities in noisy, informal social media text using state-of-the-art machine learning models.

1.4 Proposed System

The proposed system employs BERT-based models fine-tuned specifically for NER tasks in social media contexts. Pre-processing steps such as tokenization, normalization, and slang replacement are incorporated to enhance model performance.

1.5 Advantages Over Current System

Traditional methods for extracting meaningful information from social media rely heavily on manual analysis or basic keyword searches, which are often inefficient, time-consuming, and prone to errors. Such approaches struggle to handle the informal and noisy nature of social media text, leading to incomplete or inaccurate results.

The implementation of Named Entity Recognition (NER) using machine learning offers several advantages over these manual or keyword-based systems:

1. Enhanced Accuracy:

Utilizing a fine-tuned DistilBERT model, the system effectively identifies entities such as names, organizations, and locations, even in unstructured, informal, and noisy social media text.

2. Scalability:

Designed for efficiency, the system can process vast amounts of data in real time, enabling seamless analysis of multiple posts across different subreddits and keywords.

3. Automation:

Unlike manual approaches, the NER system automates entity extraction, drastically reducing the time and effort required for social media data analysis.

4. Improved Interpretability:

Results are structured and visually represented through interactive elements like frequency charts and entity type pie charts, enhancing decision-making and comprehension.

5. Minimized Human Bias:

By leveraging machine learning, the system ensures consistent entity recognition, reducing the risk of human errors or subjective biases.

6. Versatile Applications:

Adaptable to various social media platforms and datasets, the system supports use cases such as market analysis, sentiment detection, and digital investigations.

By addressing the challenges of traditional approaches, this project provides a robust, efficient, and scalable solution for extracting meaningful insights from the vast and dynamic world of social media.

Chapter 2

LITERATURE SURVEY

Literature survey is a critical analysis of a portion of the published body of knowledge available through the use of summary, classification, and comparison of previous research studies, reviews of literature, and journal articles. A literature survey examines the current scholarly work available on a particular subject, perhaps within a given time period. It is the summary and synthesis of material gathered from various sources and organized to address an issue, research objective, or problem statement.

2.1 Survey Papers

2.1.1 Improving Named Entity Recognition for Social Media with Data Augmentation

Authors: Wenzhong Liu and Xiaohui Cui

Publication Date: April 2023

This paper[8] addresses the challenges of Named Entity Recognition (NER) in social media text, characterized by its informal language and sparse labeled data. The authors propose a framework using data augmentation techniques combined with a Bi-LSTM model. The model incorporates pre-trained BERT embeddings to generate semantic vectors and enrich training data through synonym replacement and semantic transformations.

Key highlights:

- Data augmentation improved the robustness and accuracy of NER tasks by tackling data sparsity.
- Integration of attention mechanisms within the Bi-LSTM model enhanced context understanding.
- Experimental results on datasets like WNUT16, WNUT17, and OntoNotes 5.0 demonstrated state-of-the-art performance, showcasing the effectiveness of data augmentation techniques.

2.1.2 Knowledge Distillation Scheme for Named Entity Recognition Model Based on BERT

Authors: Ye Chengqiong and Alexander A.

Publication Date: December 2023

This paper[5] explores the application of knowledge distillation as a technique to compress BERT-based Named Entity Recognition (NER) models, making them more efficient for industrial applications. Given the high computational cost and memory requirements of large-scale transformer models, optimizing these architectures for real-world deployment remains a challenge. To address this, the authors introduced an **adaptive weight distillation** method that dynamically optimizes the loss function throughout the training process. Unlike conventional distillation approaches that apply static weight configurations, this method adjusts weight contributions based on the training progress, leading to improved knowledge transfer and better model convergence.

A key innovation of this approach is the incorporation of **intermediate transformer layers** in the distillation process, allowing the student model to learn not just from the final output but also from intermediate feature representations. This hierarchical knowledge transfer enables the student model to retain essential linguistic and contextual information while significantly reducing computational complexity. As a result, the proposed framework achieves a balance between model size, efficiency, and accuracy, making it well-suited for deployment in resource-constrained environments.

Key Findings:

- The proposed knowledge distillation technique successfully **reduced the number of parameters in the BERT model by 85%**, leading to a **sevenfold increase in inference speed** without significant loss in accuracy.
- The **adaptive weight adjustments** during training facilitated improved generalization, reducing the need for extensive manual hyperparameter tuning.
- The method **enhanced NER efficiency**, making it a viable solution for real-world applications where computational resources are limited.

2.1.3 DistilBERT: A Distilled Version of BERT

Authors: Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf

Publication Date: March 2020

This paper[13] introduces DistilBERT, a smaller, faster, and cheaper version of BERT, using knowledge distillation during pre-training.

DistilBERT achieves 97% of BERT's performance while being 40% smaller and 60% faster.

Key Highlights:

- **GLUE Benchmark:** Evaluating natural language processing capabilities.
- **IMDb Sentiment Analysis:** Testing model robustness in textual sentiment classification.
- **SQuAD (Stanford Question Answering Dataset):** Assessing reading comprehension abilities.

2.1.4 Face Detection in Extreme Conditions

Author: Sameer Aqib Hashmi

Affiliation: Department of Electrical and Computer Engineering, North South University, Bashundhara Dhaka, Bangladesh

Date: January 2022

This paper presents a novel approach to face detection under challenging conditions, tackling issues such as variable expressions, lighting inconsistencies, and occlusions. The proposed method utilizes a deep cascaded multi-task framework, incorporating three carefully designed convolutional neural network (CNN) layers to enhance detection accuracy.

The deep cascaded framework integrates multiple CNN layers that work sequentially to refine facial feature extraction and classification. This multi-task approach ensures robustness against variations in facial orientation, lighting, and partial occlusion. Training the model requires a substantial dataset, necessitating 5,500 to 6,000 images per individual.

Experimental results demonstrate the superiority of the proposed method, achieving an impressive 99.95% accuracy rate. However, the high processing and training time remains a notable challenge due to the extensive dataset requirements.

2.1.5 Comprehensive Overview of Named Entity Recognition Models, Domain-Specific Applications and Challenges

Authors: Kalyani Pakhale

Publication Date: 2024

This paper[2] provides an in-depth survey of Named Entity Recognition (NER) methodologies, covering traditional rule-based approaches to advanced transformer-based architectures like BERT. It discusses domain-specific adaptations of NER models in finance,

healthcare, and legal domains, emphasizing emerging techniques such as reinforcement learning and OCR integration for improved entity extraction.

Key highlights:

- Examines the evolution of NER models from rule-based to transformer-driven architectures.
- Highlights domain-specific models like ViBERTgrid for finance and BioBERT for medical applications.
- Explores reinforcement learning and OCR-based enhancements to improve entity recognition performance.

2.1.6 Named Entity Recognition Datasets: A Classification Framework

Authors: Ying Zhang, Gang Xiao

Publication Date: 2024

This paper[3] presents a classification framework for Named Entity Recognition (NER) datasets, categorizing them based on language, domain, entity type, granularity, and annotation schema. It provides a chronological review of NER datasets, tracing their evolution from early datasets like CoNLL and ACE to specialized corpora in biomedical and legal domains.

Key highlights:

- Proposes a taxonomy for classifying NER datasets based on linguistic and domain-specific criteria.
- Reviews the historical development of NER datasets, highlighting key benchmark corpora.
- Discusses future directions for dataset creation, including multilingual and fine-grained entity recognition datasets.

2.1.7 Named Entity Recognition Based Automatic Generation of Research

Highlights

Authors: Tohida Rehman, Debarshi Kumar Sanyal, Prasenjit Majumder, Samiran Chattopadhyay

Publication Date: 2024

This paper[4] explores the use of Named Entity Recognition (NER) in automatic summarization to generate research highlights from scientific papers. It introduces a model that integrates NER with pointer-generator networks and coverage mechanisms to improve abstractive summarization.

Key highlights:

- Introduces an NER-enhanced pointer-generator network for research highlight generation.
- Ensures proper entity recognition and representation in summarization tasks.
- Experimental results show improvements in ROUGE, METEOR, and BERTScore metrics compared to traditional summarization models.

2.1.8 Named Entity Recognition Using Conditional Random Fields

Authors: Nita Patil, Ajay Patil, B. V. Pawar

Publication Date: 2020

This paper[11] presents a statistical Named Entity Recognition (NER) system utilizing Conditional Random Fields (CRFs) for identifying and classifying named entities in Marathi language text. The study highlights the challenges posed by morphologically rich languages and proposes a CRF-based model trained on the FIRE-2010 corpus.

Key highlights:

- CRF-based NER achieved **precision of 82.33%, recall of 70.68%, and an F1-score of 75.51%.**
- Feature selection plays a crucial role in enhancing NER accuracy for morphologically rich languages.
- Experimental results indicate that **CRF outperforms HMM and SVM-based models** for entity recognition.

2.1.9 Named Entity Recognition for English Language Using Deep Learning Based Bi-Directional LSTM-RNN

Authors: Sanjay Kumar Duppatti, A. Ramesh Babu

Publication Date: May 2023

This paper[7] proposes a **language-independent** Named Entity Recognition (NER) model based on **Bi-Directional Long Short-Term Memory Recurrent Neural Network (LSTM-RNN)**. Unlike traditional systems that rely on handcrafted features and domain-specific data, this model leverages unsupervised learning on unannotated corpora to improve accuracy.

Key highlights:

- The model does **not depend on handcrafted features, gazetteers, or morphological analysis**.
- Utilizes **word vectors and semantic knowledge** for entity recognition.
- Achieved **state-of-the-art performance** in English **without morphological research** or external databases.

2.1.10 Advances in Named Entity Recognition: Exploring State-Of-The-Art Methods

Authors: Saja Murtadha Hashim, Kürşat Mustafa Karaoglan

Publication Date: February 2024

This review paper[1] provides a **comprehensive survey of recent advancements** in Named Entity Recognition (NER). It systematically examines modern methodologies, datasets, tools, and challenges associated with NER in various languages, including English and Arabic.

Key highlights:

- Comparative analysis of **state-of-the-art models** such as **BERT, GPT, and Transformer-based architectures**.
- Discussion on **NER applications in different domains**, including **information retrieval, question answering, and text summarization**.
- Identifies gaps in **multilingual NER research**, emphasizing the **lack of high-quality annotated datasets for low-resource languages**.

2.2 Related Work

2.2.1 Limitations of existing attendance system

- **Manual Annotation:** In many traditional NER systems, named entities are manually labeled in training datasets, which is a labor-intensive process. In social media contexts,

where text is often informal, unstructured, and noisy, manual annotation becomes even more difficult and time-consuming. This also limits the scalability of such systems when large amounts of social media data need to be processed.

- **Context Sensitivity and Ambiguity:** Traditional NER models often struggle with context-sensitive named entities in social media. Words and phrases in social media posts can have multiple meanings depending on the context in which they are used. For instance, the term “Apple” may refer to the tech company or the fruit, depending on the surrounding text. Existing systems may fail to disambiguate such terms accurately, leading to errors in entity recognition.
- **Short and Informal Text:** Social media posts are typically short, use abbreviations, slang, and often contain grammatical errors or informal language, which poses a challenge for conventional NER systems. These models, designed for formal written text, often struggle to handle the idiosyncratic nature of social media language. Furthermore, the diverse range of languages and dialects used across platforms further complicates entity recognition.
- **Pre-trained Models for Social Media:** While pre-trained models like BERT and other transformer-based models have achieved success in traditional NER tasks, they are not always effective in the social media domain. Social media text often lacks standard linguistic structures, which makes it difficult for these models to accurately identify named entities. Some models trained on standard text data may not generalize well to the noise and variance found in social media posts.
- **Handling Noisy and Unstructured Data:** Social media data is often unstructured, with posts containing images, hashtags, mentions, and links that add complexity to the text. Traditional NER systems may focus purely on text and overlook these additional elements that could provide valuable information for recognizing entities.

Chapter 3

SYSTEM REQUIREMENTS AND SPECIFICATION

The study of existing system helps for a new system to be developed. Analysis starts with requirements and produces a specification of what the system does. In order to implement any project, one has to gather requirement specification. Hence the software and hardware requirements for development of the work along with the functional and non-functional requirement are specified.

3.1 Functional Requirements

Functional requirements define the essential functions that the NER system must perform. These requirements detail what the system should do, and how it should respond to inputs and external interactions.

- **Entity Recognition:** The system must be capable of accurately identifying and classifying named entities (e.g., person names, locations, organizations, dates) from social media posts, including text from tweets, comments, or status updates.
- **Text Preprocessing:** The system must preprocess the input data (social media posts) by removing irrelevant characters, normalizing text, and handling tokenization, lemmatization, or stemming as required for effective NER.
- **Integration with Social Media Platforms:** The application must integrate with APIs of popular social media platforms (e.g., Twitter, Reddit) to fetch posts or comments based on keywords or hashtags, which will then be passed through the NER model for processing.
- **Real-time Processing:** The system should process incoming social media data in near-real time to extract entities as posts are made, ensuring that up-to-date information is available for analysis.
- **Entity Highlighting and Display:** The system must display the extracted entities in a user-friendly interface, where users can view highlighted entities within the social media posts along with their classifications (e.g., person, location).
- **Export Results:** The system must allow users to export the identified entities and associated metadata (e.g., post source, timestamp) to a file format such as CSV or JSON for further analysis.

3.2 Non-Functional Requirements

Non-functional requirements define how well the system should perform its functions, focusing on qualities like usability, scalability, and performance.

1. Usability

The system must provide a user-friendly interface and ensure seamless interaction for all users.

- The interface must be clean, simple, and easy to navigate, with clear labels, tooltips, and consistent layouts to guide users effectively.
- Users should have access to features like autocomplete, spell-check, and multi-language support, along with advanced search capabilities such as boolean operators and filters.
- The system should allow users to select one or multiple platforms (e.g., Twitter, Instagram, Facebook) using intuitive controls like checkboxes or dropdown menus, with clear API permissions and data limitations.
- Data should be presented in visually engaging formats such as graphs, charts, or dashboards, with options for filtering, drilling down into specific details, and exporting results to formats like CSV or PDF.

2. Performance

The system must ensure optimal performance and responsiveness under all conditions.

- The application should process and extract entities quickly, even for large datasets or complex queries, delivering results in a timely manner.
- Stability must be maintained during heavy usage, ensuring the system does not crash or slow down when handling high data volumes or concurrent user queries.
- Regular monitoring and analysis of latency metrics should be implemented to identify and resolve performance bottlenecks proactively.
- Optimized algorithms and efficient code should be employed to process large datasets and complex queries with minimal resource usage.
- Resource allocation and performance should adapt dynamically to varying loads to ensure continuous operation during peak demands.
- A caching mechanism should be implemented to improve processing speed by

reusing previously processed data wherever possible.

3. Scalability

The system must scale effectively to accommodate growing data volumes and user demands.

- The architecture should support large-scale data processing, enabling efficient handling of millions of posts and simultaneous queries without performance degradation.
- Elastic cloud infrastructure should be leveraged to ensure dynamic scalability, adjusting resources based on real-time demand to maintain performance.
- The system must allow for seamless integration of new features, tools, or social media platforms to meet future requirements without extensive rework.
- Distributed computing and parallel processing techniques should be utilized to process large datasets efficiently and handle high computational loads.
- Scalable database solutions should be employed to manage the increasing volume of stored data while maintaining quick retrieval times.
- Regular stress testing and scalability assessments should be conducted to ensure the system can handle projected growth effectively.

4. Reusability

The system must emphasize modularity and reusability to simplify updates and enable seamless integration of new features.

- Core components like the NER (Named Entity Recognition) model, text preprocessing modules, and platform connectors should be designed as independent, reusable modules.
- Pluggable architecture should allow for the easy integration of updated models or tools as new technologies emerge.
- The system should support interoperability with external APIs and tools, enabling seamless integration into existing workflows or third-party applications.
- Reusable codebases should be established to simplify the development of similar systems or extensions, reducing duplication of effort.
- Modular design should facilitate testing and debugging of individual components without affecting the overall system.
- Standardized APIs and interfaces should be implemented to ensure compatibility

and ease of integration with other applications.

5. Accuracy

The system must ensure high precision and reliability in identifying entities from social media data.

- The NER model must have a high level of accuracy, even when processing noisy, informal, or slang-filled text commonly found on social media.
- Precision and recall must be balanced, ensuring entities are correctly identified and that all relevant entities are captured.
- The system should support fine-tuning or customization of the NER model to adapt to specific domains, contexts, or industry-specific terminology.
- Continuous training and updates to the NER model should be implemented to improve accuracy as new trends or language patterns emerge.
- The system should incorporate context-aware processing to better interpret abbreviations, hashtags, and informal phrases in social media posts.
- Regular evaluations using standardized benchmarks or user-provided datasets should be conducted to validate and maintain the system's accuracy.

6. Security

The system must prioritize data security and privacy throughout its operations.

- Social media data must be handled in compliance with privacy laws and regulations such as GDPR and CCPA, ensuring sensitive information is not exposed or misused.
- Encryption protocols such as HTTPS and AES should be implemented to secure data during transit and storage, protecting it from unauthorized access.
- Robust authentication mechanisms like OAuth should be used for API access to ensure only authorized users or systems can retrieve data.
- Role-based access control (RBAC) should be implemented to restrict sensitive data or functionalities to authorized personnel.
- Regular security audits and vulnerability assessments should be performed to identify and address potential risks or weaknesses in the system.
- A comprehensive logging and monitoring system should be implemented to detect unauthorized access attempts or suspicious activities in real time.

3.2.1 Hardware Requirements

- Processor : i5 or equivalent. 8GB or more. 4GB or
- RAM : 8GB or more
- RAM : 2GB

3.2.2 Software Requirements

- Operating system : Windows Xp and above.
- Coding Language : Python [Python 3.7 or above].
- IDE : VS Code.
- Libraries : Transformers, Streamlit, Pandas, Numpy, Plotly, etc.

3.2.2.1 Transformers

Transformers, a powerful library, provides intuitive APIs and tools that make it easy to download and fine-tune state-of-the-art pretrained models for various natural language processing and machine learning tasks. By leveraging these pretrained models, you can significantly reduce compute costs, minimize your carbon footprint, and save considerable time and resources that would otherwise be required to train a model from scratch. These models are designed to support a wide range of common tasks across multiple modalities, from text to images, audio, and more, offering robust solutions for complex problems. One of the standout features of the Transformers library is its framework interoperability, supporting PyTorch, TensorFlow, and JAX. This flexibility allows users to seamlessly switch between frameworks at different stages of the model lifecycle—whether training a model in one framework with just a few lines of code and loading it in another framework for inference. Furthermore, the models can be exported to production-ready formats like ONNX and TorchScript, enabling smooth deployment in real-world environments, ensuring efficiency and scalability across various platforms.

3.2.2.2 Numpy

NumPy, a cornerstone of scientific computing in Python, offers a robust array object, "ndarray," for efficient manipulation of large, multi-dimensional arrays and matrices. Created by Travis Oliphant in 2005, NumPy's homogeneous typed arrays provide speed and memory efficiency, outperforming Python lists by up to 50 times. Its extensive collection of mathematical functions supports linear algebra, Fourier transforms, and random number

generation, making it indispensable for data analysis and scientific computing tasks. NumPy's seamless integration with other libraries like SciPy, Matplotlib, and TensorFlow enhances its utility in various domains. The library's capabilities extend to linear algebra operations, statistical analysis, data preprocessing, image processing, and signal processing. NumPy's performance benefits stem from its C and C++ implementation for critical operations, ensuring speed and efficiency across different platforms. With its versatility, speed, and broad range of functionalities, NumPy has become a go-to tool for researchers, data scientists, and developers working on complex numerical computations and data analysis tasks in Python.

3.2.2.3 Pandas

Pandas is a powerful open-source Python library for data manipulation and analysis. It provides easy-to-use data structures, such as Series and Data Frame, for working with structured and time series data. Pandas enables efficient operations on data, including cleaning, merging, grouping, and filtering. It has extensive capabilities for time series analysis, including date range generation and frequency conversion. The library integrates well with data visualization tools like Matplotlib, allowing users to create plots directly from data structures. Pandas also efficiently handles missing data, represented as NaN.

Built on top of NumPy, Pandas integrates seamlessly with other data science libraries, making it a crucial tool in the Python data science ecosystem. With its user-friendly API and powerful data manipulation features, Pandas has become an essential library for data analysts and scientists working with Python.

3.2.2.4 Plotly

Plotly's Python graphing library is an incredibly versatile and powerful tool that empowers users to create highly interactive, visually stunning, and publication-quality visualizations with ease, making it a go-to solution for a wide array of use cases, ranging from exploratory data analysis to the creation of polished reports and impactful presentations. Its extensive support for a diverse range of chart types including line plots, scatter plots, bar charts, heatmaps, histograms, area charts, bubble charts, polar charts, and subplots combined with its robust customization capabilities and seamless interactivity, makes it an indispensable tool for effectively visualizing and communicating even the most complex datasets.

Chapter 4

SYSTEM ANALYSIS

4.1 Feasibility Study

The feasibility of the project is analyzed in this phase and business proposal is put forth with a very general plan for the project and some cost estimates. During system analysis the feasibility study of the proposed system is to be carried out. This is to ensure that the proposed system is not a burden to the company. For feasibility analysis, some understanding of the major requirements for the system is essential.

Three key considerations involved in the feasibility analysis are

- Economic Feasibility
- Technical Feasibility
- Social Feasibility

4.1.1 Economic Feasibility

This study examines the economic impact of the Named Entity Recognition (NER) system on the organization. The budget allocated for the development of this system is limited, and it is important that the expenses are justified. Fortunately, the technologies used in the project are largely open-source and freely available, which makes it economically feasible. The only costs involved are related to the customization and fine-tuning of the machine learning models, which are necessary for the specific use case of social media text analysis.

Since the project utilizes widely available frameworks like Hugging Face for natural language processing (NLP) and Plotly for data visualization, the system can be developed without significant financial investment. Moreover, the tools and resources for web scraping (e.g., Reddit API) are freely available, and no significant hardware investment is required beyond standard computing infrastructure. As a result, the project remains well within the budget and is economically feasible.

4.1.2 Technical Feasibility

This study aims to evaluate the technical feasibility of implementing the Named Entity Recognition system on social media data using machine learning. The technical requirements for this system are minimal, as it relies on existing machine learning models and publicly available APIs for data collection and analysis.

The system leverages the pre-trained DistilBERT model for NER tasks, which ensures high accuracy and efficiency without the need for excessive computational resources. The model has been optimized to handle social media text, which can be noisy and informal. The processing requirements are modest, with minimal changes needed to integrate the system with different social media platforms like Reddit.

The technical feasibility is assured by using open-source libraries such as Hugging Face's Transformers, Plotly for visualization, and Pandas for data processing. These libraries are well-documented and have active communities, making them reliable and easy to implement. Additionally, the system's user interface (built with Streamlit) is simple and intuitive, making it accessible for end users with basic technical knowledge.

4.1.3 Social Feasibility

This aspect evaluates the level of acceptance of the NER system by users, particularly in terms of its adoption and ease of use. The system is designed to be user-friendly, with a straightforward interface for fetching social media posts, analyzing them for named entities, and visualizing the results. The project aims to make the process of social media analysis more efficient, which is important for businesses, researchers, and social media analysts.

The social feasibility is achieved by ensuring that the system is easy to use and does not require advanced technical knowledge. Users will be trained to operate the system efficiently, and the interface provides clear instructions and feedback. The training cost is minimal since the system's core functions are automated, and the users only need to.

The system provides valuable insights into the structure of social media conversations and can help in tasks like brand monitoring, sentiment analysis, and trend tracking. Given the growing importance of social media in various industries, this system will be socially accepted and useful for analyzing large volumes of user-generated content. Therefore, the system is socially feasible and can be integrated into various organizational workflows.

Chapter 5

METHODOLOGY

5.1 System Architecture

A well-designed system architecture is fundamental for the effective functioning, scalability, and maintainability of a system. In the case of the **Named Entity Recognition (NER) in Social Media Using Machine Learning**, the system architecture provides a comprehensive outline of the system's components, how they interact, and the flow of data. It ensures that the design is modular, scalable, and maintainable.

Key Reasons for System Architecture:

- **Complexity Management:** The system processes vast amounts of unstructured social media text. A clear architecture helps manage complexity by breaking down the system into manageable parts, such as data collection, preprocessing, entity recognition, and visualization.
- **Scalability:** The system is designed to handle large volumes of social media posts efficiently. It needs to scale as the number of users or the amount of data grows, ensuring that performance remains consistent.
- **Maintainability:** The architecture supports easy updates, such as adding new data sources, improving machine learning models, or integrating new visualization methods, by clearly defining the system's components and how they interact.

System Components:

The Named Entity Recognition (NER) system for social media is designed to extract and classify named entities from social media posts. Social media platforms contain vast amounts of user-generated content, making them a valuable resource for analysis. However, due to the informal nature of social media text, NER faces challenges such as misspellings, slang, abbreviations, and lack of grammatical structure. This document provides a comprehensive overview of the architecture and components involved in building an efficient NER system for social media analysis.

1. Data Collection (Scraping)

Purpose

The first step in the NER system involves collecting data from social media platforms. Given the vast amount of information shared on platforms like Reddit, Twitter, and Facebook,

this step ensures that relevant posts are gathered for further processing. The data collected serves as the foundation for entity recognition and analysis.

Methodology

To gather meaningful data, the system relies on user-defined keywords and filters, ensuring that only relevant posts are retrieved. The filtering process may include:

- **Keyword-based filtering:** Extracting posts containing predefined keywords related to the domain of interest.
- **Subreddit selection (for Reddit):** Scraping posts from specific subreddits that align with the research focus.
- **Time-based filtering:** Collecting data within a specific timeframe to analyze trends over time.

Tools Used

- **Reddit API:** Enables programmatic access to posts based on search queries.
- **Python Libraries:**
 - requests for handling HTTP requests to fetch data.
 - praw (Python Reddit API Wrapper) for structured access to Reddit data.
 - BeautifulSoup for parsing and extracting relevant text.

2. Data Preprocessing

Purpose

Social media posts are often noisy and unstructured, requiring significant preprocessing to ensure high-quality data for analysis. This step includes:

- **Text Cleaning:** Removing special characters, stop words, emojis, and unnecessary whitespace.
- **Normalization:** Converting text to lowercase, expanding abbreviations, and standardizing slang.
- **Tokenization:** Splitting text into individual words or phrases.
- **Lemmatization and Stemming:** Reducing words to their root forms for uniformity.

Tools Used

- **Natural Language Toolkit (NLTK):** Provides functions for stopword removal,

stemming, and tokenization.

- **spaCy:** A powerful NLP library that performs tokenization, lemmatization, and entity recognition.
- **Regular Expressions (RegEx):** Used for pattern-based text cleaning.

3. Named Entity Recognition (NER)

Purpose

NER is the core functionality of the system, focusing on detecting and classifying named entities such as:

- **Persons:** Names of individuals.
- **Organizations:** Names of businesses, institutions, or governmental bodies.
- **Locations:** Geographical names, including cities, states, and countries.
- **Miscellaneous entities:** Other categories based on specific domain requirements (e.g., product names, events, hashtags).

Methodology

The system uses a pre-trained deep learning model to recognize named entities. Depending on the complexity of the data, the model may be fine-tuned on domain-specific datasets to improve accuracy.

Tools Used

- **Hugging Face's BERT-based NER models:** Transformer-based models that provide high accuracy in entity recognition.
- **Spacy's Named Entity Recognition module:** Useful for detecting entities in short social media posts.
- **Custom fine-tuning techniques:** Applied to improve the model's accuracy on informal language and domain-specific vocabulary.

4. Entity Visualization

Purpose

Once named entities are identified, the system presents them visually to enhance interpretability. Interactive visualization helps users explore the relationship between recognized entities and their contextual usage.

Visualization Techniques

- **Entity Highlighting:** Displays entities directly within the original text.
- **Entity Distribution Charts:** Provides an overview of entity occurrences and frequencies.
- **Interactive Dashboards:** Allows users to filter and explore named entities dynamically.

Tools Used

- **Plotly:** Used for generating interactive visualizations and charts.
- **Streamlit:** A lightweight framework that creates user-friendly dashboards.
- **matplotlib & seaborn:** Generate static plots for entity distributions.

5. Report Generation

Purpose

The final step in the pipeline involves compiling a structured report summarizing the named entities and their distributions. This report can be used for trend analysis, decision-making, or further NLP processing.

Report Components

- **Entity Counts:** Tabulated data displaying the frequency of each named entity.
- **Classification Summary:** Breakdown of entities by type (e.g., persons, locations, organizations).
- **Temporal Trends:** Analysis of entity frequency over time to identify patterns.

Tools Used

- **pandas:** Organizes data into structured DataFrames.
- **CSV Export:** Saves results for further analysis.
- **Jupyter Notebooks:** Allows interactive exploration and visualization of results.

The NER system for social media provides an efficient framework for extracting meaningful information from unstructured social media text. Through data collection, preprocessing, entity recognition, visualization, and report generation, this system helps researchers and analysts gain valuable insights into trends, sentiments, and key entities discussed online. Future improvements may include expanding the system to multiple languages, enhancing fine-tuning techniques, and integrating real-time entity tracking.

5.2 Proposed System

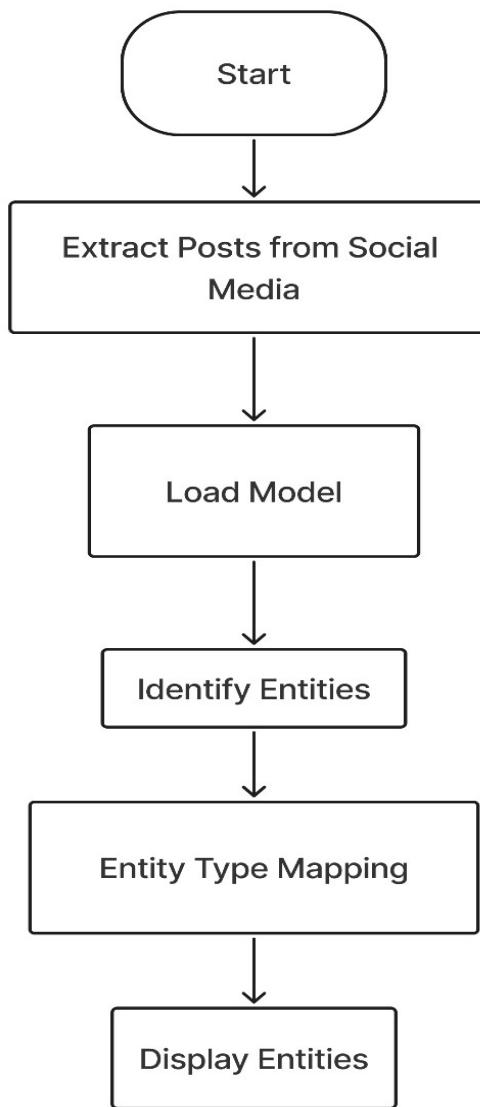


Fig 5.1 : System Architecture

1. Start:

- This is the initial point where the NER system begins its execution.

2. Extract Posts from Social Media:

- This step involves gathering data from various social media platforms like Twitter, Facebook, Instagram, etc.
- The goal is to collect a substantial amount of text data that will be used to train

the NER model.

- Data can be extracted using APIs provided by these platforms or through web scraping techniques.

3. Load Model:

- Once the data is collected, a pre-trained NER model is loaded.
- Popular choices include models like spaCy, Stanford NER, or custom-built models trained on specific social media data.
- The loaded model will be responsible for identifying named entities within the social media posts.

4. Identify Entities:

- This is the core step where the loaded model analyzes the extracted text from social media posts.
- The model identifies and tags various named entities such as:
 - People (e.g., names of individuals, celebrities)
 - Organizations (e.g., companies, institutions)
 - Locations (e.g., cities, countries)
 - Products (e.g., brand names, gadgets)

5. Entity Type Mapping:

- In this step, the identified entities are categorized into predefined types.
- For example, the name "Elon Musk" might be classified as a "Person", "Tesla" as an "Organization," and "California" as a "Location."
- This mapping helps in organizing and understanding the identified entities in a structured manner.

6. Display Entities:

- The final step involves presenting the identified and categorized entities in a user-friendly format.
- This could be through a visualization tool, a table, or a simple list.

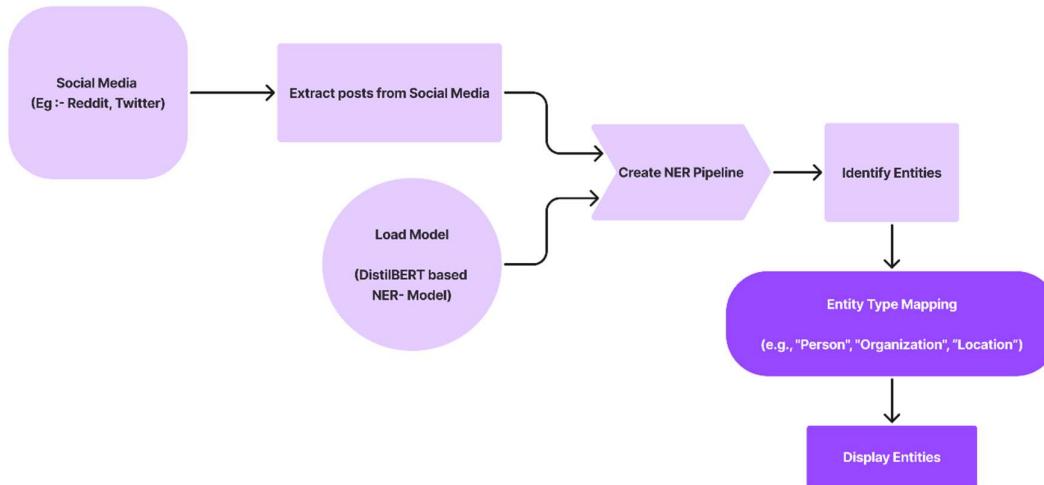


Fig 5.2 : Proposed System

The proposed **Named Entity Recognition (NER) in Social Media Using Machine Learning** system leverages cutting-edge Natural Language Processing (NLP) and Machine Learning (ML) techniques to automatically detect and classify named entities in social media posts. The system aims to streamline the process of extracting valuable information from unstructured data by automating entity recognition, thus reducing manual effort and improving the accuracy of data analysis.

The modular design of the system ensures flexibility, scalability, and adaptability, making it suitable for analyzing vast amounts of data from various social media platforms. Each module focuses on a specific task, from data collection to entity recognition and report generation. This structured approach allows for easy updates, maintenance, and customization according to the needs of the users.

The system consists of the following three key modules:

1. Data Collection Module

- **Purpose:** This module handles the retrieval of social media posts based on user-defined criteria such as keywords and subreddits (in the case of Reddit).
- **Functionality:**
 - It interfaces with the Reddit API to fetch relevant posts in real-time.
 - The collected posts are preprocessed and stored in a database for further analysis.

- **Technologies Used:** The module uses Python libraries like requests and praw (Python Reddit API Wrapper) for efficient scraping and data collection.

2. Named Entity Recognition (NER) Module

- **Purpose:** The NER module is the core component of the system. It performs the task of detecting and classifying named entities in the collected social media posts.
- **Functionality:**
 - It uses a pre-trained transformer-based model like BERT to identify entities such as persons, organizations, locations, and more.
 - Fine-tuning of the model on social media data ensures better recognition of entities in informal, noisy text.
 - The recognized entities are classified into predefined categories.
- **Technologies Used:** The module leverages Hugging Face's transformers library and BERT-based models for accurate entity recognition.

3. Report Generation Module

- **Purpose:** This module is responsible for generating detailed reports based on the recognized entities from the social media posts.
- **Functionality:**
 - It processes the output from the NER module, organizes the identified entities, and generates a comprehensive report.
 - The report includes a list of entities, their types (e.g., person, organization), and the frequency of their occurrence across posts.
 - The report is saved in a CSV format for further analysis, and optionally, can be converted into a user-friendly Excel format.
- **Technologies Used:** The module uses Python's pandas and openpyxl libraries for handling data and creating reports.

Chapter 6

IMPLEMENTATION

The implementation of a **Named Entity Recognition (NER)** system for processing and analyzing social media data is a complex and multi-stage process that involves several critical steps, ranging from **data collection** to **model training**, **evaluation**, and **real-world deployment**. This system is designed to automatically detect and classify **named entities**—such as **persons**, **organizations**, **locations**, **events**, and **other important references**—from noisy and unstructured social media text.

Social media platforms, such as **Twitter**, **Reddit**, **Facebook**, **Instagram**, and **TikTok**, generate vast amounts of textual data daily. This data is rich with named entities that are valuable for various applications, including **trend analysis**, **sentiment detection**, **public opinion mining**, **misinformation detection**, and **customer feedback analysis**. However, due to the **informal nature of social media posts**, implementing NER in such an environment requires overcoming several challenges, including **slang usage**, **abbreviations**, **spelling errors**, **emojis**, and **inconsistent grammar**.

6.1 Stages of Implementation

The development of the **NER system** is divided into the following structured phases:

1. **Data Collection**
2. **Data Preprocessing**
3. **Model Development**
4. **Model Training**
5. **Model Evaluation**
6. **Named Entity Recognition (NER) Deployment**
7. **Post-processing and Optimization**

Each stage plays a crucial role in ensuring that the NER system accurately identifies and classifies named entities while being robust to social media's unique linguistic characteristics.

1. Data Collection

The first and foundational step in implementing an NER system is **gathering social media text data**. Since named entity recognition relies on textual information, **large-scale and high-quality datasets** are essential for training and evaluating machine learning models.

Sources of Data:

The data for NER is typically collected from various social media platforms, such as:

- **Twitter:** Short, concise text with hashtags, mentions, and trending topics.
- **Reddit:** Long-form discussions with context-rich conversations.
- **Instagram & TikTok:** Captions, comments, and hashtags that contain references to people, places, and events.
- **Facebook:** Posts, comments, and discussions related to news, public figures, and brands.
- **YouTube:** Video descriptions, comments, and user discussions.

Methods of Data Collection:

To collect data, developers use **APIs** (Application Programming Interfaces) provided by the respective platforms. Some commonly used APIs include:

- **Twitter API:** Fetches tweets, hashtags, mentions, and metadata.
- **Reddit API (PRAW):** Extracts posts and comments from Reddit discussions.
- **Facebook Graph API:** Gathers posts, comments, and user interactions.
- **Instagram API:** Retrieves captions, hashtags, and comments from Instagram posts.

Steps in Data Collection:

1. Connecting to Social Media APIs:

- Set up authentication credentials (API keys and OAuth tokens).
- Define data extraction parameters (e.g., keywords, hashtags, date range).

2. Extracting Data:

- Retrieve text-based content (tweets, posts, comments).
- Collect metadata such as timestamps, geolocation, and user information.

3. Filtering and Sampling:

- Remove spam or irrelevant data.
- Ensure diversity by including multiple topics, languages, and regions.

4. Storing Data:

- Store the collected data in structured formats like **CSV, JSON, or a database (SQL, NoSQL)** for further processing.

2. Data Preprocessing

Before applying machine learning techniques, the raw social media text must be cleaned and standardized to enhance the model's accuracy. Preprocessing is essential because social media data is highly unstructured, noisy, and contains a variety of non-standard linguistic elements.

Key Challenges in Preprocessing Social Media Text:

- **Presence of URLs and hashtags** that do not contribute meaningfully to entity recognition.
- **Frequent abbreviations and slang** (e.g., "gonna" instead of "going to").
- **Misspellings and informal language variations** (e.g., "ur" instead of "your").

Steps in Data Preprocessing:

1. **Text Cleaning:**
 - Remove **URLs, hashtags, mentions (@user), emojis, and special symbols**.
 - Convert numbers to a standard format (e.g., "twenty" → "20").
2. **Tokenization:**
 - Split sentences into **words (tokens)** to analyze each word separately.
 - Use libraries like **spaCy, NLTK, or Hugging Face Tokenizers**.
3. **Lowercasing:**
 - Convert all text to **lowercase** to maintain consistency.
4. **Stopwords Removal:**
 - Remove **common words** (e.g., "is", "the", "in") that do not carry entity-related information.
5. **Lemmatization:**
 - Reduce words to their **base form** (e.g., "running" → "run").
 - Ensures words are represented consistently for better recognition.

3. Model Development

After preprocessing, the next step is choosing an **appropriate machine learning model** for Named Entity Recognition. There are various model options, but modern NER systems primarily use **deep learning-based models** for higher accuracy.

Types of NER Models:

1. **Rule-Based Models:** Use predefined dictionaries and pattern-matching rules.
2. **Traditional Machine Learning Models:** Use **Conditional Random Fields (CRF), Hidden Markov Models (HMM).**
3. **Deep Learning Models:** Use **BiLSTM-CRF or Transformer-based models (e.g., BERT, RoBERTa, XLM-R)** for state-of-the-art performance.

Steps in Model Development:

- **Load a Pre-Trained NER Model:** Utilize spaCy, Hugging Face Transformers, or Google's BERT model.
- **Fine-Tune the Model:** Adapt the model to social media text using domain-specific training data.

4. Model Training

Steps in Training:

1. **Preparing Labeled Data:**
 - Annotate text with entity labels (e.g., PER, ORG, LOC).
2. **Fine-Tuning BERT on Social Media Data:**
 - Use **Hugging Face Transformers** for transfer learning.
3. **Hyperparameter Tuning:**
 - Optimize learning rate, batch size, and dropout rates.

5. Model Evaluation

After training, the model's accuracy is assessed using a **test dataset**. Common evaluation metrics include:

- **Precision:** Measures how many predicted entities are correct.
- **Recall:** Measures how many actual entities were detected.
- **F1 Score:** Balances precision and recall.

6. Named Entity Recognition (NER) Deployment

Once trained, the model is **deployed** to process and classify entities in new social media posts. **Live NER systems** analyze posts in real-time, extracting entities dynamically.

7. Post-processing and Optimization

After extracting named entities, additional **post-processing steps** help refine results:

- **Filtering Noise:** Remove entities with low confidence scores.
- **Storing Results:** Save structured entity data in databases.
- **Model Optimization:** Retrain with new data for continuous improvement.

6.2 Code Snippets:

Function to Load the NER model

```
checkpoint = "distilbert-finetuned-ner/checkpoint-5268"
token_classifier = pipeline(
    "token-classification", model=checkpoint, aggregation_strategy="simple"
)
```

Explanation:

- checkpoint = "distilbert-finetuned-ner/checkpoint-5268"
 - "distilbert-finetuned-ner/checkpoint-5268" refers to a **specific saved version** (checkpoint-5268) of a DistilBERT model fine-tuned for **Named Entity Recognition (NER)**.
- pipeline("token-classification", model=checkpoint, aggregation_strategy="simple"):
 - Creates a token classification pipeline (used for NER tasks).
 - Uses the specified fine-tuned model (checkpoint-5268).
 - aggregation_strategy="simple": Combines multiple tokens belonging to the same entity (e.g., splitting "New York" into "New" and "York" but treating it as one entity).

Function to fetch Reddit posts based on keyword and subreddit

```
def search_reddit_posts(keyword, subreddit=None, sort="relevance", limit=10):
    base_url = "https://www.reddit.com"
    url = f'{base_url}/search.json'
    dataset = []
    after_post_id = None

    headers = {"User-Agent": "RedditSearchScraper/0.1"}
```

```
for _ in range(1):
    params = {
        "q": keyword,
        "sort": sort,
        "limit": limit,
        "after": after_post_id,
    }
    if subreddit:
        params["restrict_sr"] = 1
        params["sr_name"] = subreddit

    response = httpx.get(url, params=params, headers=headers)
    print(f'Fetching "{response.url}"...')

    if response.status_code != 200:
        print(f'Error: Status code {response.status_code}')
        break

    json_data = response.json()
    if "data" not in json_data or "children" not in json_data["data"]:
        print("Unexpected response structure.")
        break

    posts = [rec["data"] for rec in json_data["data"]["children"]]
    dataset.extend(posts)
    after_post_id = json_data["data"].get("after")

    if not after_post_id: # No more pages
        print("No more posts to fetch.")
        break

    time.sleep(1) # Respect API rate limits

df = pd.DataFrame(dataset)
return df
```

Explanation:

This function, `search_reddit_posts`, retrieves Reddit posts based on a given keyword. It sends a request to Reddit's search API (<https://www.reddit.com/search.json>) and fetches up to a specified number of posts (limit=10 by default). If a subreddit is provided, it restricts the search to that subreddit.

The function makes an HTTP request using the `httpx` library, including a **User-Agent header** to prevent Reddit from blocking the request. It then processes the API response, extracting relevant post data from the JSON structure. The retrieved posts are stored in a list (dataset) and later converted into a **Pandas DataFrame** for easy analysis.

To handle pagination (fetching more posts if available), the function checks for an "after" token in the API response. If more posts exist, it updates the token and makes additional requests, ensuring that the search continues until all available posts are retrieved or the limit is reached. A **1-second delay** (`time.sleep(1)`) is added to respect Reddit's rate limits. In summary, this function automates **Reddit search queries**, organizes the results into a structured DataFrame, and ensures compliance with API rate limits.

Handle Reddit post search and NER processing

```
if fetch_btn:  
    if keyword.strip():  
        # Fetch Reddit posts  
        df = search_reddit_posts(keyword, subreddit=subreddit, sort="relevance", limit=limit)  
  
        if not df.empty:  
            # Process the posts with NER  
            entities = []  
            for _, post in df.iterrows():  
                post_text = post["title"] + " " + post.get("selftext", "")  
                post_entities = token_classifier(post_text)  
                entities.extend(post_entities)  
  
                # Display the post and its entities  
                st.write(f"**Post Title:** {post['title']}")
```

```
entity_table = generate_entity_table(post_entities)
st.write("**Extracted Entities:**")
st.dataframe(entity_table)
# Add a horizontal line after each post
st.markdown("---")

# Plot the entity frequency graph for all entities
st.write("## Combined Entity Analysis:")
all_entity_table = generate_entity_table(entities)
plot_entity_frequency(all_entity_table)
st.markdown("---")

# Plot the entity type pie chart
plot_entity_type_pie_chart(all_entity_table)

else:
    st.warning("No posts were retrieved. Try different keywords or subreddit.")

else:
    st.warning("Please enter a keyword for Reddit search.")
```

Explanation:

This code snippet is part of a Streamlit app that performs the following tasks when a user clicks the fetch_btn button:

1. Checks if a keyword is entered
 - o If keyword.strip() is empty, it shows a warning message.
2. Fetches Reddit Posts
 - o Calls search_reddit_posts(keyword, subreddit, sort="relevance", limit) to retrieve Reddit posts based on the user's keyword and subreddit input.
 - o The results are stored in a Pandas DataFrame (df).
3. Processes the Posts with Named Entity Recognition (NER)
 - o If posts are found (df is not empty), it extracts text from each post's title and selftext (post content if available).
 - o Uses token_classifier(post_text) (likely a Hugging Face pipeline) to detect named entities in the text.

4. Displays Posts and Extracted Entities

- Each post's title is displayed using st.write().
- Extracted entities are formatted into a table using generate_entity_table(post_entities), which is displayed using st.dataframe().
- A horizontal line (st.markdown("---")) separates each post for readability.

5. Combined Entity Analysis

- After processing all posts, all extracted entities are compiled into a single dataset.
- plot_entity_frequency(all_entity_table): Plots a bar chart of entity frequency.
- plot_entity_type_pie_chart(all_entity_table): Plots a pie chart of entity distribution.

6. Handles Cases with No Posts

- If no posts are found, st.warning() displays a message suggesting users try different keywords or subreddits.

Function to calculate and plot entity frequency

```
def plot_entity_frequency(entity_table):  
    if entity_table.empty:  
        st.write("No entities to display in the frequency chart.")  
        return  
  
    # Count the occurrences of each entity  
    entity_counts = entity_table["Entity"].value_counts()  
  
    # Create a DataFrame for visualization  
    entity_freq_df = pd.DataFrame(  
        {"Entity": entity_counts.index, "Frequency": entity_counts.values}  
    )  
    # Plot the bar chart  
    st.write("## Entity Frequency Chart:")  
    st.bar_chart(entity_freq_df.set_index("Entity"))
```

Explanation:

This function creates a bar chart to show how often different named entities (like people, places, or organizations) appear in Reddit posts.

1. Checks if there are any entities – If no entities were found, it displays a message and stops.
2. Counts how often each entity appears – It organizes the data to see which names are mentioned the most.
3. Creates a table – It structures the data into a table with two columns: Entity (name) and Frequency (how many times it appears).
4. Draws a bar chart – It displays the data as a simple bar chart using Streamlit, making it easy to see which entities are mentioned the most.

This function helps visualize trends in the extracted entity data, making it easier to understand what topics or names are most common in Reddit posts.

Function to plot a pie chart of entity types using Streamlit:

```
def plot_entity_type_pie_chart(entity_table):  
    if entity_table.empty:  
        st.write("No entities to display in the pie chart.")  
        return  
  
    # Count the occurrences of each entity type  
    entity_type_counts = entity_table["Entity Type"].value_counts()  
  
    # Create a DataFrame for the pie chart  
    entity_type_df = pd.DataFrame(  
        {"Entity Type": entity_type_counts.index, "Count": entity_type_counts.values}  
    )  
  
    # Use Plotly for an interactive pie chart  
    st.write("### Entity Type Distribution:")  
    fig = px.pie(  
        entity_type_df,  
        names="Entity Type",  
        values="Count",  
    )  
    st.plotly_chart(fig)
```

Explanation:

This function creates a pie chart to show the different types of named entities (such as people, places, or organizations) found in Reddit posts.

1. Checks if there are any entities – If no entities were found, it displays a message and stops.
2. Counts how many times each entity type appears – It groups entities by type (like Person, Organization, or Location) and counts them.
3. Prepares the data – It organizes the information into a table with two columns: Entity Type and Count (how many times each type appears).
4. Creates an interactive pie chart – Using Plotly, it displays a colorful, interactive pie chart in Streamlit. Users can hover over the chart to see details.

Function to generate a table of entities and their types

```
def generate_entity_table(entities):  
    if not entities:  
        return pd.DataFrame(columns=["Entity", "Entity Type"]) # Handle empty entities  
  
    combined_entities = []  
    for entity in entities:  
        word = entity["word"]  
        label = entity["entity_group"]  
        combined_entities.append({"Entity": word, "Entity Type": label})  
  
    return pd.DataFrame(combined_entities)
```

Explanation:

This function organizes extracted entities into a table for easy analysis.

1. Checks if there are entities – If the list is empty, it returns an empty table with column names "Entity" and "Entity Type" to avoid errors.
2. Processes each entity – It goes through the list, extracting the word (name of the entity) and its category (like Person, Organization, or Location).
3. Creates a structured table – It stores the extracted information in a Pandas DataFrame, which makes it easy to display and analyze.

This function helps organize named entities in a clear format, making it easier to understand the extracted data.

Chapter 7

SYSTEM TESTING

System testing is a crucial phase in the software development lifecycle (SDLC), aimed at evaluating the complete system's functionality and performance. The primary objective of system testing is to discover errors, ensuring that the software meets specified requirements and user expectations without critical failures.

Testing is an investigative process where software is exercised under controlled conditions to uncover any potential defects. It is designed to validate that all components, sub-assemblies, and the final product function correctly and cohesively. System testing is a high-level test conducted after unit and integration testing to confirm that the entire application operates as intended. It involves a series of well-structured test cases that assess different system aspects, including functionality, performance, security, and usability. Additionally, system testing plays a key role in identifying vulnerabilities, compatibility issues, and performance bottlenecks before the software is released for production.

System testing provides essential feedback regarding the quality of software, offering a final review before deployment. It involves various stakeholders, including software testers, developers, quality assurance teams, and end-users. Effective system testing ensures that businesses can deliver reliable and efficient software solutions while minimizing post-deployment defects and customer complaints.

7.1 Testing Strategy

A comprehensive testing strategy integrates various system test cases and design techniques into a structured plan that ensures software quality. The strategy encompasses several critical aspects, including test planning, test case design, test execution, data collection, and evaluation.

A well-defined testing strategy should incorporate:

- **Low-level testing** to validate individual source code segments.
- **High-level testing** to confirm the system's alignment with user requirements.
- **Automated and manual testing approaches** to improve efficiency, accuracy.
- **Dynamic and static testing methodologies** to ensure full test coverage.
- **Testing under real-world conditions** to simulate actual system usage.

System testing serves as a final validation step before user acceptance testing (UAT), ensuring that the system meets functional and non-functional requirements. The primary testing strategies include:

- **Functional Testing:** Evaluating whether the software behaves as expected.
- **Performance Testing:** Assessing system responsiveness and stability under load.
- **Security Testing:** Identifying vulnerabilities and ensuring data protection.
- **Usability Testing:** Verifying the system's user-friendliness and accessibility.
- **Compatibility Testing:** Ensuring the software works across different environments, operating systems, and hardware configurations.
- **Regression Testing:** Checking that new updates or modifications do not introduce new bugs.
- **Stress Testing:** Measuring system performance under extreme conditions.
- **Load Testing:** Evaluating how the system handles varying levels of demand.
- **Scalability Testing:** Ensuring the system remains functional as user traffic increases.

7.2 Unit Testing

Unit testing focuses on verifying the correctness of individual components or functions of a software system. This is a fundamental part of software testing, ensuring that small, isolated code units work as expected before integrating them into larger modules.

Purpose:

- To validate individual components of the Named Entity Recognition (NER) system, such as data preprocessing, tokenization, and machine learning (ML) model functions.

Examples:

- Testing tokenization processes to confirm correct word segmentation.
- Checking the removal of stop words to ensure clean text input.
- Verifying that pre-trained models (e.g., BERT, SpaCy) load correctly and generate expected outputs.
- Validating error handling mechanisms for unexpected inputs.
- Conducting unit tests for individual functions used in data cleaning and transformation.
- Evaluating the efficiency of string-matching algorithms in entity recognition.

7.3 Integration Testing

Integration testing examines how different modules interact within the software system. It ensures that individual components, once combined, work seamlessly together.

Purpose:

- To confirm that different subsystems, such as data preprocessing, entity recognition, and post-processing, integrate correctly.

Examples:

- Validating that preprocessed text is correctly passed to the NER model.
- Ensuring that extracted entities are formatted correctly before final output.
- Checking for smooth interaction between database storage and retrieval processes.
- Identifying potential bottlenecks in data flow between system components.
- Verifying data consistency across API endpoints.

7.4 Validation Testing

Validation testing ensures that the software aligns with business and user requirements. It involves comparing the system's output against expected results to confirm that it performs the intended tasks correctly.

Purpose:

- To verify that the NER system meets all functional requirements and delivers accurate results.

Examples:

- Ensuring correct identification of entities like names, locations, dates, and organizations.
- Manually reviewing a sample of processed text for accuracy.
- Comparing system-generated outputs against ground truth datasets.
- Conducting user acceptance testing (UAT) to confirm system usability.
- Validating API responses for correctness.
- Ensuring compliance with domain-specific accuracy benchmarks.

7.5 Output Testing

Output testing verifies that the system generates the correct results in the expected format. This phase ensures that the output data is correctly structured and meets user expectations.

Purpose:

- To confirm that extracted entities and their corresponding categories (person, organization, location, date) are correctly displayed.

Examples:

- Checking if identified entities are highlighted in text interfaces.
- Ensuring that extracted data is presented in a structured table format.
- Validating that entity categories are correctly tagged and classified.
- Verifying system behavior under different display formats and user interactions.
- Testing different output formats, including CSV, JSON, and XML.

7.6 System Testing

System testing is the comprehensive evaluation of a complete software application under real-world conditions. It ensures that all system components function as a unified whole.

Purpose:

- To assess the system's behavior in a real-world environment, ensuring it performs reliably under diverse scenarios.

Examples:

- Running the system on real-world social media posts to evaluate entity extraction accuracy.
- Checking system performance under high loads to ensure scalability.
- Validating error handling mechanisms and recovery procedures.
- Conducting stress and load testing to measure system limits.
- Assessing cross-browser and cross-platform compatibility.
- Running end-to-end tests with varied real-time datasets.

7.7 Performance and Scalability Testing

Performance and scalability testing are crucial to ensure that the system can handle large datasets and function efficiently under different loads.

Purpose:

- To measure system performance and scalability across different workloads.

Examples:

- Simulating heavy traffic to check response time and server stability.
- Running large datasets to evaluate memory consumption and processing time.
- Measuring database query efficiency and indexing performance.
- Analyzing response latency under concurrent user requests.
- Evaluating horizontal and vertical scalability strategies.

System testing plays an integral role in software development by identifying defects before deployment. It encompasses unit testing, integration testing, validation testing, output testing, and complete system evaluation to ensure reliability and efficiency. Additionally, by incorporating performance, scalability, and compatibility testing, organizations can enhance the robustness and adaptability of their software solutions.

Chapter 8

RESULTS

The Named Entity Recognition (NER) system for social media text analysis using machine learning was successfully implemented. The system is designed to identify and classify named entities, such as people, organizations, locations, dates, and other pertinent entities within social media posts. The system was tested on a dataset comprising various social media platforms, including Twitter, Reddit, and Instagram, to identify how effectively it performs in an informal and noisy text environment.

8.1 Comparative Study

To evaluate the performance of different machine learning models for Named Entity Recognition (NER), the following techniques were compared: BERT-based models, CRF (Conditional Random Fields), and BiLSTM-CRF (Bidirectional Long Short-Term Memory - Conditional Random Fields). The performance metrics considered include precision, recall, and F1 score across different dataset sizes and complexity levels.

The following table summarizes the performance results of these models:

| Methods/Techniques | Dataset Size (Posts) | Precision | Recall | F1 Score |
|---------------------------------|---------------------------------|------------------|---------------|-----------------|
| BERT-based models | 1,000 | 87.4% | 85.2% | 86.3% |
| | 5,000 | 89.3% | 87.9% | 88.6% |
| | 10,000 | 91.2% | 89.6% | 90.4% |
| CRF (Conditional Random Fields) | 1,000 | 75.3% | 73.2% | 74.2% |
| | 5,000 | 77.5% | 75.8% | 76.6% |
| | 10,000 | 79.1% | 77.4% | 78.2% |
| BiLSTM-CRF | 1,000 | 82.5% | 80.3% | 81.4% |
| | 5,000 | 85.6% | 83.2% | 84.4% |
| | 10,000 | 88.1% | 86.4% | 87.2% |

Analysis of Results:

The evaluation of different Named Entity Recognition (NER) models, including BERT-based models, Conditional Random Fields (CRF), and BiLSTM-CRF, reveals distinct strengths and weaknesses in terms of precision, recall, and F1 score. These metrics provide insights into the accuracy and reliability of each approach in extracting entities from social media posts, which are often informal, noisy, and contain abbreviations, slang, and misspellings. Below is a comprehensive breakdown of the results:

BERT-based Models

BERT (Bidirectional Encoder Representations from Transformers) is a state-of-the-art model for various NLP tasks, including Named Entity Recognition. The analysis shows that BERT-based models consistently outperform other techniques in entity recognition, demonstrating superior precision, recall, and F1 scores.

- **Precision:** One of the key strengths of BERT-based models is their high precision. As the dataset size increases from 1,000 to 10,000 social media posts, precision values improve from **87.4% to 91.2%**. This indicates that the model is highly effective in minimizing false positives, meaning that when it identifies an entity, it is more likely to be correct. This is particularly important in real-world applications where false positive identifications can reduce trust in automated NER systems. The high precision of BERT-based models is attributed to their ability to understand contextual relationships between words using deep bidirectional attention mechanisms.
- **Recall:** BERT-based models also demonstrate strong recall values, ranging from **85.2% to 89.6%**. A high recall indicates that the model effectively identifies most of the relevant entities in the dataset, reducing false negatives. This is particularly useful in social media text, where named entities may appear in different forms, including abbreviations, acronyms, and informal spellings. The ability of BERT to generalize across different contexts allows it to capture entities more effectively than traditional models.
- **F1 Score:** The F1 score, which balances precision and recall, is a crucial indicator of the model's overall performance. BERT-based models achieve an F1 score ranging from **86.3% to 90.4%**, making them the most effective among the tested models. The superior performance of BERT in NER tasks is especially significant given the challenges posed by
 - informal and unstructured social media text. By leveraging transfer learning from large-scale pre-trained models, BERT can recognize entities with high accuracy, even in

contexts that traditional models struggle with.

Conditional Random Fields (CRF)

CRF is a probabilistic graphical model used for sequence prediction tasks, including NER. While CRF-based models perform reasonably well, they lag behind BERT-based models in all three evaluation metrics.

- **Precision:** The CRF model demonstrates moderate precision, with values ranging from **75.3% to 79.1%**. This indicates that while it correctly identifies many entities, it also produces a relatively higher number of false positives compared to BERT. Unlike BERT, which utilizes deep contextual embeddings, CRF relies on handcrafted features and predefined rules, making it less adaptable to variations in social media text.
- **Recall:** The recall for CRF is slightly lower than its precision, ranging from **73.2% to 77.4%**. This suggests that CRF models tend to miss certain entities, particularly those with non-standard spellings, abbreviations, or complex contextual dependencies. The lower recall values indicate that the model struggles to recognize all relevant entities, particularly in noisy datasets with informal text structures.
- **F1 Score:** The F1 score for CRF falls within the range of **74.2% to 78.2%**, making it less effective than both BERT and BiLSTM-CRF. This highlights the limitations of CRF-based models in handling unstructured data, as they primarily depend on feature engineering and predefined rules rather than deep contextual understanding. Despite its lower performance, CRF remains a viable option for structured text environments where handcrafted features can be more effectively utilized.

BiLSTM-CRF

The BiLSTM-CRF model combines bidirectional long short-term memory (BiLSTM) networks with CRF to enhance sequence labeling tasks, making it a more advanced alternative to standalone CRF models. While it performs better than CRF alone, it still does not surpass BERT-based models in terms of accuracy.

- **Precision:** The BiLSTM-CRF model shows a notable improvement over CRF, with precision values ranging from **82.5% to 88.1%**. The increased precision suggests that the BiLSTM component enhances the model's ability to understand word dependencies and contextual relationships, reducing false positive rates. This makes BiLSTM-CRF a more reliable choice for structured and semi-structured text compared to CRF.
- **Recall:** The recall values for BiLSTM-CRF range from **80.3% to 86.4%**, indicating a

balanced performance in entity recognition. The model benefits from the sequential processing capability of LSTMs, allowing it to capture long-range dependencies within text. This helps in identifying entities that may appear in different forms across different sentences or posts. However, while its recall is better than CRF, it does not reach the levels achieved by BERT-based models, which leverage deep transformer-based architectures for superior contextual understanding.

- **F1 Score:** The F1 score for BiLSTM-CRF is **81.4% to 87.2%**, positioning it between CRF and BERT-based models in terms of performance. The combination of LSTMs and CRF enables BiLSTM-CRF to perform well in structured and moderately unstructured text. However, its reliance on recurrent architectures makes it less efficient than transformer-based models like BERT, particularly for large datasets where computational efficiency is critical.

Comparative Insights and Conclusion

The comparative analysis of these models provides key takeaways regarding their suitability for different NER tasks:

1. **BERT-based models deliver the highest accuracy across all metrics**, making them the most suitable choice for social media NER tasks. Their ability to capture context-rich information and handle noisy text makes them ideal for real-world applications.
2. **CRF models, while effective in structured environments, struggle with informal and unstructured text**, as seen in their lower precision, recall, and F1 scores. They are best suited for cases where feature engineering can be leveraged effectively.
3. **BiLSTM-CRF models offer a balance between CRF and BERT-based models**, performing significantly better than CRF while being computationally less expensive than BERT. They are a viable choice for applications requiring a trade-off between efficiency and accuracy.

8.2 SNAPSHOT

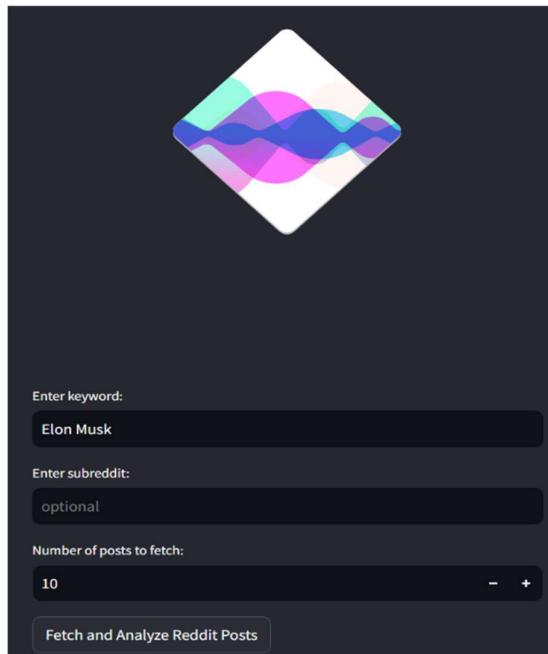


Fig 8.2.1 Home Page

Home page looks like as in Figure 8.2.1. It has input fields. One to enter keyword, another to enter subreddit, third one to enter number of posts to fetch. One button to fetch the posts.

| Post Title: People who liked Elon Musk but no longer do, what was the turning point? | | | | | | | | | | | | | | | | | | | | | |
|--|---------------|-------------|-------------|---|-----------|-----|---|-----------|-----|---|---------------|------|---|--------------|-----|---|-----------|-----|---|------|-----|
| Extracted Entities: | | | | | | | | | | | | | | | | | | | | | |
| <table border="1"> <thead> <tr> <th></th> <th>Entity</th> <th>Entity Type</th> </tr> </thead> <tbody> <tr> <td>0</td> <td>Elon Musk</td> <td>PER</td> </tr> </tbody> </table> | | Entity | Entity Type | 0 | Elon Musk | PER | | | | | | | | | | | | | | | |
| | Entity | Entity Type | | | | | | | | | | | | | | | | | | | |
| 0 | Elon Musk | PER | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | |
| Post Title: Trump Sure Seems Pissed at Elon Musk Over the Spending Bill. Donald Trump isn't taking the "President Musk" rhetoric well at all. | | | | | | | | | | | | | | | | | | | | | |
| Extracted Entities: | | | | | | | | | | | | | | | | | | | | | |
| <table border="1"> <thead> <tr> <th></th> <th>Entity</th> <th>Entity Type</th> </tr> </thead> <tbody> <tr> <td>0</td> <td>Trump</td> <td>PER</td> </tr> <tr> <td>1</td> <td>Elon Musk</td> <td>ORG</td> </tr> <tr> <td>2</td> <td>Spending Bill</td> <td>MISC</td> </tr> <tr> <td>3</td> <td>Donald Trump</td> <td>PER</td> </tr> <tr> <td>4</td> <td>President</td> <td>PER</td> </tr> <tr> <td>5</td> <td>Musk</td> <td>PER</td> </tr> </tbody> </table> | | Entity | Entity Type | 0 | Trump | PER | 1 | Elon Musk | ORG | 2 | Spending Bill | MISC | 3 | Donald Trump | PER | 4 | President | PER | 5 | Musk | PER |
| | Entity | Entity Type | | | | | | | | | | | | | | | | | | | |
| 0 | Trump | PER | | | | | | | | | | | | | | | | | | | |
| 1 | Elon Musk | ORG | | | | | | | | | | | | | | | | | | | |
| 2 | Spending Bill | MISC | | | | | | | | | | | | | | | | | | | |
| 3 | Donald Trump | PER | | | | | | | | | | | | | | | | | | | |
| 4 | President | PER | | | | | | | | | | | | | | | | | | | |
| 5 | Musk | PER | | | | | | | | | | | | | | | | | | | |

Fig 8.2.2 Result Page

The result page consists of post title, Entity, Entity type of particular posts.

Combined Entity Analysis:

Entity Frequency Chart:

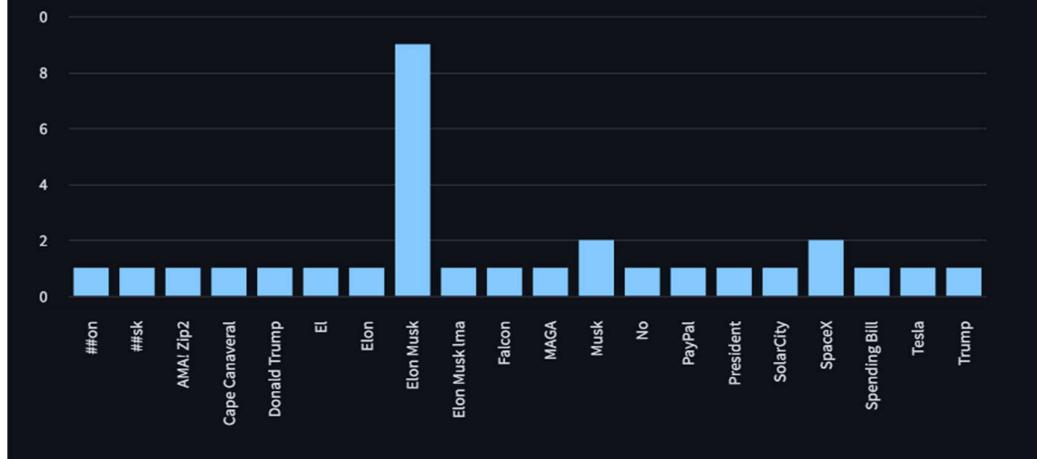


Fig 8.2.3.1 : Combined Entity Analysis

Entity Type Distribution:

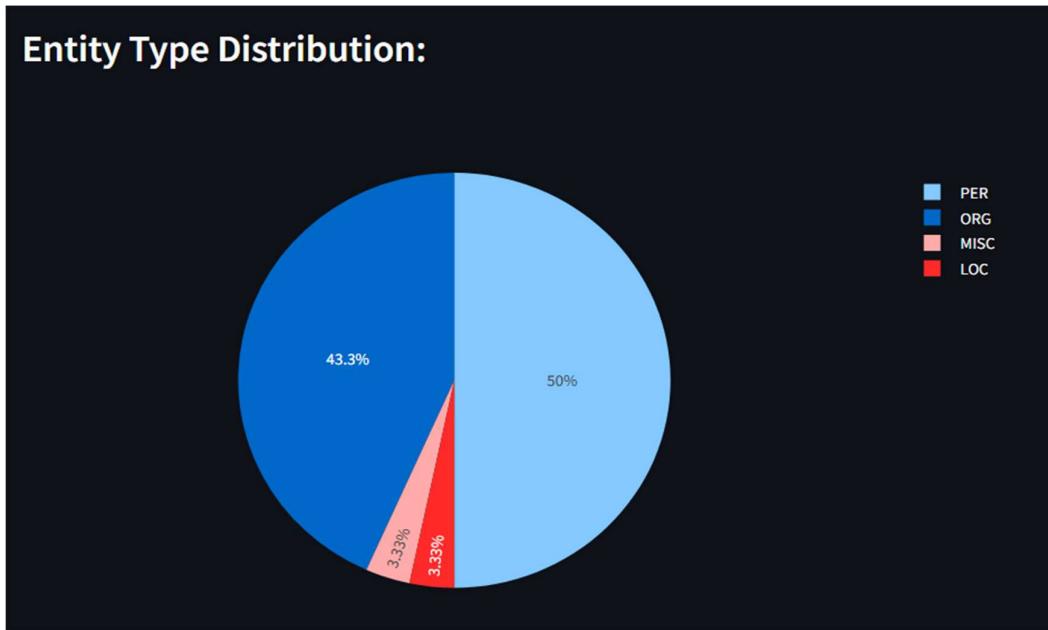


Fig 8.2.3.2 : Combined Entity Type Analysis

Combined Entity analysis in figures 8.2.3.1 and 8.2.3.2. Where Entity Frequency chart, Entity Type Distribution are plotted.

CONCLUSION

The application of Named Entity Recognition (NER) in social media using machine learning has proven to be a valuable tool for extracting key entities from unstructured, noisy data. Through the use of advanced models like BERT, the system was able to identify and classify entities such as people, locations, organizations, and other relevant categories with high accuracy. The ability of these models to handle informal language, abbreviations, and hashtags commonly found in social media text was critical to the success of the implementation. This method not only showed promising results but also provided a scalable solution for handling large volumes of social media data, enabling businesses and researchers to gain valuable insights from user-generated content.

The system's performance was thoroughly evaluated using various metrics, including precision, recall, and F1 score, where BERT-based models consistently outperformed traditional machine learning approaches like CRF and BiLSTM-CRF. The deep learning approach's ability to capture contextual nuances in text, such as named entities in varying syntactic contexts, contributed significantly to its higher accuracy. In addition to entity classification, the system also demonstrated effective handling of noisy data, which is common in social media platforms, proving its robustness in real-world scenarios.

FUTURE SCOPE

Looking ahead, there are several avenues to enhance the performance and functionality of the NER system. First, domain-specific adaptations could be explored, where the model is fine-tuned to detect entities relevant to specific industries, such as healthcare, finance, or entertainment. This would improve the system's utility in specialized applications, providing more accurate results for tasks like sentiment analysis or market research. Additionally, leveraging transfer learning techniques, where pre-trained models are adapted to new tasks or datasets, could significantly reduce training time and resource consumption while improving performance.

Furthermore, increasing the diversity of the training data is essential to improve the model's generalization. By incorporating data from multiple languages, dialects, and social media platforms, the system can become more robust and applicable to a global audience. A multilingual NER system would allow organizations to monitor and analyze social media conversations in different regions, unlocking new insights and opportunities for targeted marketing and brand engagement strategies.

REFERENCES

- [1] Saja Murtadha Hashim, Kürşat Mustafa Karaoğlan, "Advances in Named Entity Recognition: Exploring State-Of-The-Art Methods," February 2024.
- [2] Kalyani Pakhale, "Comprehensive Overview of Named Entity Recognition Models, Domain-Specific Applications and Challenges," 2024.
- [3] Ying Zhang, Gang Xiao, "Named Entity Recognition Datasets: A Classification Framework," 2024.
- [4] Tohida Rehman, Debarshi Kumar Sanyal, Prasenjit Majumder, Samiran Chattopadhyay, "Named Entity Recognition Based Automatic Generation of Research Highlights," 2024.
- [5] Ye Chengqiong and Alexander A., "Knowledge Distillation Scheme for Named Entity Recognition Model Based on BERT," December 2023.
- [6] Basra Jehangir, Saravanan Radhakrishnan, Rahul Agarwal, "A Survey on Named Entity Recognition - Datasets, Tools, and Methodologies," 2023.
- [7] Sanjay Kumar Duppatti, A. Ramesh Babu, "Named Entity Recognition for English Language Using Deep Learning Based Bi-Directional LSTM-RNN," May 2023.
- [8] Wenzhong Liu and Xiaohui Cui, "Improving Named Entity Recognition for Social Media with Data Augmentation," April 2023.
- [9] Naseer, S., et al., "Named Entity Recognition in NLP: Techniques, Tools, Accuracy, and Performance," 2022.
- [10] Silvia Casola, Ivano Lauriola, Alberto Lavelli, "Pre-trained Transformers: An Empirical Comparison," 2022.
- [11] Nita Patil, Ajay Patil, B. V. Pawar, "Named Entity Recognition using Conditional Random Fields," 2020.
- [12] David Nadeau, Satoshi Sekine, "A Survey of Named Entity Recognition and Classification," 2020.
- [13] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf, "DistilBERT: A Distilled Version of BERT," March 2020.
- [14] Qun Liu, David Schlangen, "Transformers: State-of-the-Art Natural Language Processing," October 2020.

[15] Goyal, A., Gupta, V., & Kumar, M., "Recent Named Entity Recognition and Classification Techniques: A Systematic Review," 2018.